

关于“智能评阅算法效果的多维度综合评价模型构建及应用”的研究

摘要

在新质生产力发展背景下，人工智能驱动的考生评卷智能化成为教育数字化改革的重要方向。针对主流评卷系统主观题人工评阅效率瓶颈，某实验室提出“一人工+双 AI”协同评卷机制（两类 AI 算法背对背评分并与人工一评交叉验证）。本研究基于附件 1 中填空与简答题的人工及双 AI 评阅数据，围绕以下核心问题展开：

针对问题一：通过描述性统计分析与分布拟合检验模型，定量刻画了三种评阅方式的集中趋势、离散程度、偏态和峰态特性。首先，计算了样本均值、方差和标准差，以评估数据的基本特性。接着，通过偏度与峰度指标，分析数据的对称性及尖峭程度。最终，进行正态性检验，比较三类评阅方式的数据分布特征，为智能评卷的应用提供了科学依据。该研究为理解人工智能在教育评估中的作用奠定了基础。

针对问题二：通过定义准确性和稳定性指标，包括 Pearson 相关系数、平均绝对误差和评分波动性等，量化评阅结果的可靠性。运用 AHP 方法确定各指标权重，并通过 TOPSIS 法计算相对贴近度，实现对算法效果的综合评价。研究结果将为智能评阅技术的优化与应用提供科学依据，推动教育评估的数字化转型。

针对问题三：在问题二的基础上，针对不同学科的智能评阅效果进行了跨学科的评价对比。建立了多级评价体系，通过动态权重调整引入学科难度系数，提升评估的可比性。采用鲁棒性设计处理全零数据和极端值，确保评价的稳定性。结合 AHP 方法确定权重，并利用 ICC 和 SI 指标验证算法的稳定性。研究结果为不同学科的评阅效果提供了综合评价，助力智能评阅技术的优化与应用。

针对问题四：通过研究提出了一种基于动态分配与人工复核的智能评阅方案，结合两类人工智能算法以优化评分流程。针对大分值低容忍误差的主观题，优先使用准确性高的算法 A，并进行人工复核；而对小分值高容忍误差的客观题，则使用效率高的算法 B，进行抽查。设计了误差阈值判定和算法分配策略，确保评分的准确性与公平性。同时，通过历史误差数据的分析，不断优化分配规则。该方案实现了评分自动化，并通过数据驱动提升了整体评阅效率。

关键词：描述性统计分析 分布拟合检验模型 TOPSIS-加权模型 多级评价体系 动态权重调整 鲁棒性设计 可视化验证方法 AI-算法动态分配与人工复核方案

（一）问题重述

1.1 问题概述和明确方向

2024 年初中国提出加快发展新质生产力，人工智能作为培育新动能的核心领域，其在考试评卷智能化中的应用成为人才选拔领域的重要突破点。当前主流网上评卷系统虽实现客观题自动批改，但主观题仍依赖人工多评，存在人力消耗大的问题；而某实验室基于大模型、手写识别、自然语言处理等技术，提出“一人工+双 AI”的协同机制（即两种智能算法背对背评分并与人工一评结果交叉验证），推动智能评阅从理论探索迈向规模化应用。附件 1 提供了两类人工智能算法与人工评阅的成绩数据（含填空题型和简答题型，不同子对象满分值各异）。

在此背景下，需明确以下方向：

一是分析人工评阅与两类 AI 算法评阅数据的分布特征；

二是从多维度构建评价指标体系并设计综合评价模型以评估算法效果；

三是基于附件 2 的其他科目数据开展学科维度的评阅效果对比；

四是结合附件 3 的科目题号、分值及误差阈值设计 AI 算法使用方案并进行可行性分析。

1.2 问题设计流程图

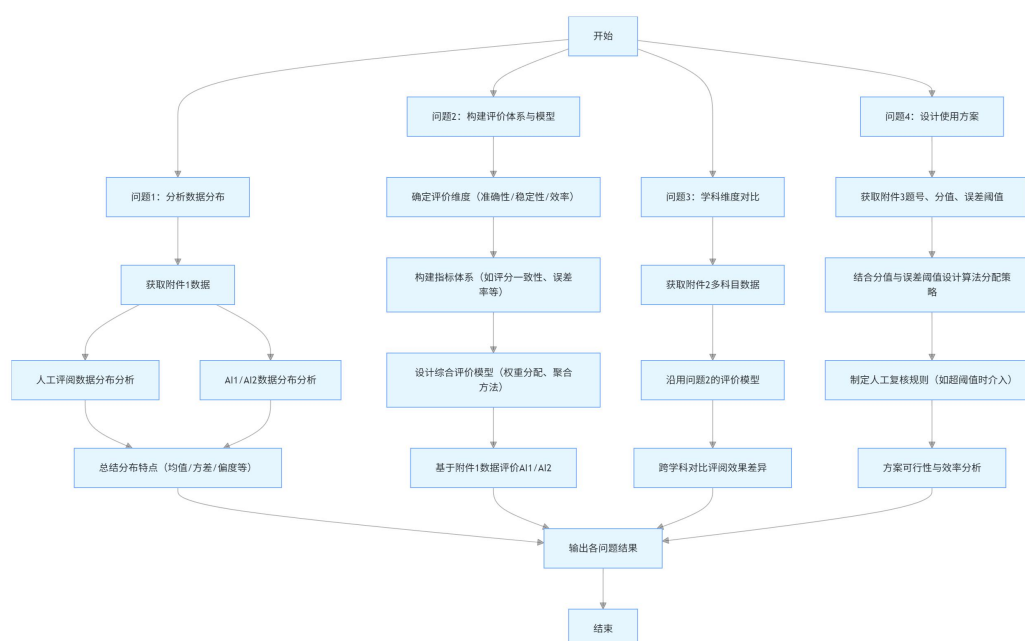


图 1 问题设计流程图

(二) 符号说明

符号	含义
s^2	样本方差
s	样本标准差
Skewness	偏度
Kurtosis	峰度
n	被评阅的样本总数
m	评价指标总数
\hat{y}_i	智能算法对第 <i>i</i> 个样本的评分
y_i	人工对第 <i>i</i> 个样本的最终评分
S	题目满分值
$n_k (n_k \geq 1)$	第 <i>k</i> 类题目样本数
$S_k (S_k \geq 0)$	第 <i>k</i> 类题目满分
$w_j (\sum w_j = 1)$	第 <i>j</i> 项指标权重
$Ci [0,1]$	算法综合评价得分
n	总题数
W_i	第 <i>i</i> 题权重
V_i	第 <i>i</i> 题分值
Δ_i	$\Delta_i = S_i^{AI} - S_i^{\text{人工}} $, 第 <i>i</i> 题 AI 评分误差
S_i^{AI}	第 <i>i</i> 题 AI 评分
$S_i^{\text{人工}}$	第 <i>i</i> 题人工评分
τ_i	第 <i>i</i> 题误差阈值
T_1	分值阈值
T_2	误差阈值

（三）模型假设

1. 假设人工评阅与 AI 评阅数据的分布形态（如正态性、偏度、峰度）可通过描述性统计量（均值、方差、分位数等）准确刻画，且异常值已被合理处理，不影响整体分布特征分析。

2. 构建评价指标体系时，假设各底层指标（如评分一致性、误差率、效率指标等）之间相互独立，不存在显著的相关性，以保证综合评价模型的权重分配和指标聚合逻辑有效。

3. 假设综合评价模型中各指标的权重确定方法（如层次分析法、熵权法等）满足数学公理体系，权重系数能客观反映指标对智能评阅算法效果的重要程度。

4. 针对学科维度评价时，假设不同科目数据在模型中可通过标准化处理消除分值差异影响，且学科特性（如题型复杂度、评分规则）可通过量化指标（如误差阈值敏感度）统一度量。

5. 设计 AI 算法使用方案时，假设附件 3 中的误差阈值与科目分值匹配逻辑符合教育考试业务需求，且算法部署成本、人工复核流程的时间效率满足实际应用场景的约束条件。

（四）模型建立与求解

4.1 问题一：分析人工评阅和两类人工智能算法评阅数据的分布特征

4.1.1 模型一方法：分布特征分析模型（描述性统计分析+分布拟合与检验模型）

为全面、科学地分析三类评阅方式的数据分布特征，建议采用描述性统计分析+分布拟合与检验的建模方法。这是统计学中最优且标准的分布特征分析方法，能定量刻画数据的集中、离散、偏态、峰态等特性，并判断数据分布类型^[1]。

4.1.2 过程

设样本数据为 x_1, x_2, \dots, x_n ，本文用如下描述性统计分析+分布拟合与检验模型来实现：

首先，计算均值 \bar{x} 为：

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad (1)$$

其次，计算样本方差 s^2 与样本标准差 s ：

样本方差 s^2 ：

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (2)$$

样本标准差 s 为:

$$s = \sqrt{s^2} \quad (3)$$

于是为衡量数据分布的对称性，有偏度公式:

$$\text{Skewness} = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s} \right)^3 = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3}{s^3} \quad (4)$$

衡量数据分布的尖峭或平坦程度，峰度公式为:

$$\text{Kurtosis} = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s} \right)^4 = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^4}{s^4} \quad (5)$$

4.1.3 结果：图表、数值分析

最后进行正态性检验，判断分布类型并对比三类评阅方式数据的分布特征如下:

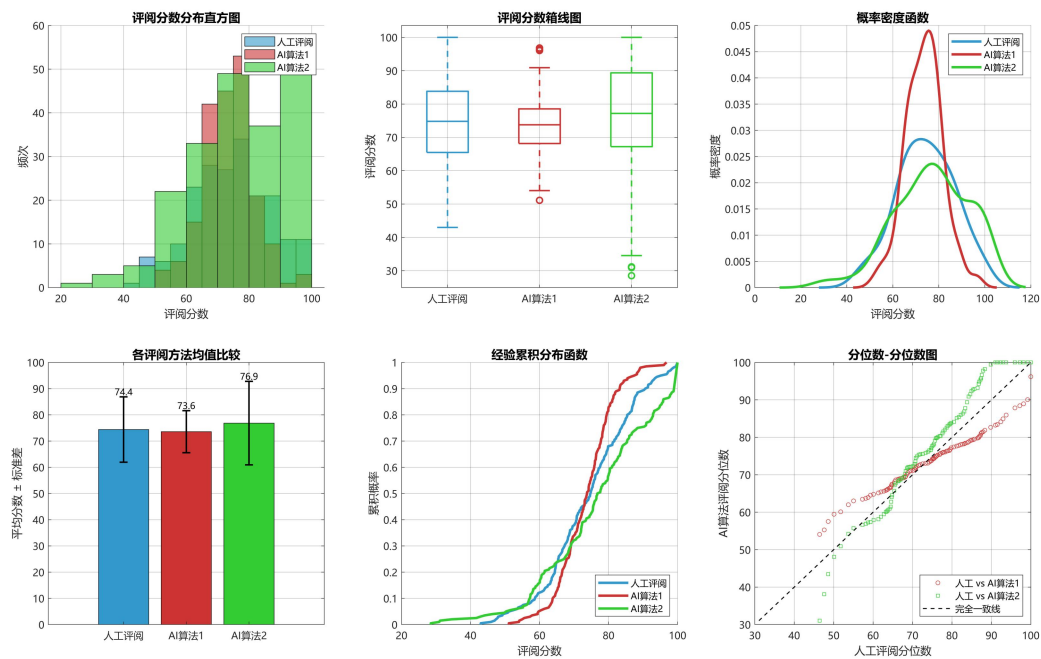


图2 三类评阅方法统计特征对比

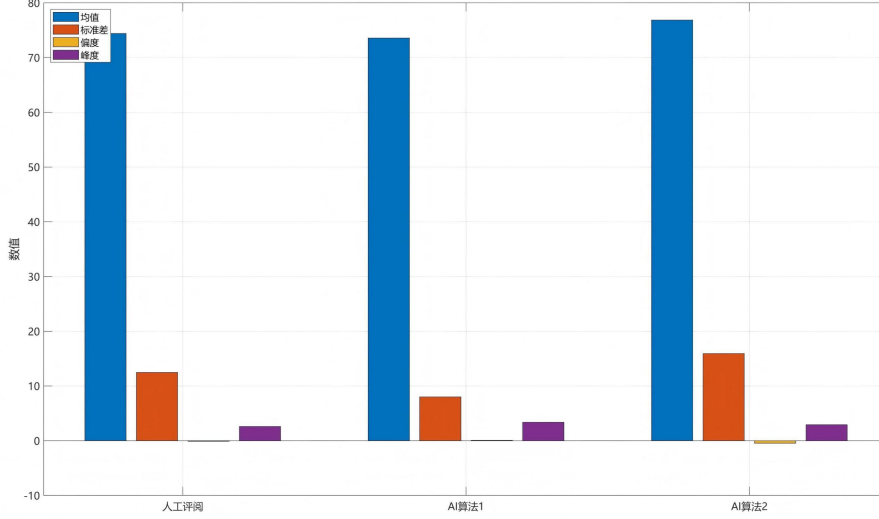


图3 人工评阅与AI算法评阅分布特点分析

4.2 问题二：选择不同评价角度，构建“智能评阅算法”的评价指标体系，设计智能评阅算法的综合评价模型，利用所给数据对两类人工智能算法给出评价

4.2.1 模型二方法：多维度综合评价模型（准确性+稳定性+TOPSIS 等加权模型）

本研究构建智能评阅算法多维度综合评价模型，定义基础参数后，以 Pearson 相关系数、MAE 等衡量准确性，评分波动性、EP 指标评估稳定性，结合 AHP 确定权重并通过一致性检验，利用 TOPSIS 法计算相对贴近度实现综合评价，为算法性能量化提供科学框架^[1]。

4.2.2 过程

设 n 为被评阅的样本总数， m 为评价指标总数，智能算法对第 i 个样本的评分， y_i 为人工对第 i 个样本的最终评分， S 为题目满分值， k 为题目类型编号（ $k \in \{1, 2, 3, 4\}$ 对应 T11, T13, T14, T15），接下来将会从评价指标计算、综合评价模型两部分展开：

评价指标计算分为准确性指标和稳定性指标，首先是准确性指标：

有 Pearson 相关系数如下：

$$r_k = \frac{\sum_{i=1}^n (y_i^{(k)} - \bar{y}^{(k)}) (\hat{y}_i^{(k)} - \bar{\hat{y}}^{(k)})}{\sqrt{\sum_{i=1}^n (y_i^{(k)} - \bar{y}^{(k)})^2 \sum_{i=1}^n (\hat{y}_i^{(k)} - \bar{\hat{y}}^{(k)})^2}} \quad (6)$$

其次是平均绝对误差：

$$MAE_k = \frac{1}{n_k} \sum_{i=1}^{n_k} |y_i^{(k)} - \hat{y}_i^{(k)}| \quad (7)$$

计算均方根误差如下：

$$\text{RMSE}_k = \sqrt{\frac{1}{n_k} \sum_{i=1}^{n_k} (y_i^{(k)} - \hat{y}_i^{(k)})^2} \quad (8)$$

准确性指标中最后一项完全一致率计算如下：

$$\text{CR}_k = \frac{1}{n_k} \sum_{i=1}^{n_k} \Pi(y_i^{(k)} = \hat{y}_i^{(k)}) \quad (9)$$

此外是稳定性指标，首先是评分波动性：

$$\sigma_k = \sqrt{\left[(1/n_k) \sum (\hat{y}_i^{(k)} - y_i^{(k)} - \mu_k)^2 \right]} \quad (10)$$

$$\mu_k = \frac{1}{n_k} \sum_{i=1}^{n_k} (\hat{y}_i^{(k)} - y_i^{(k)})^2 \quad (11)$$

接着是极端偏差比例：

$$\text{EP}_k = \frac{1}{n_k} \sum_{i=1}^{n_k} \Pi(|\hat{y}_i^{(k)} - y_i^{(k)}| > 0.2S_k) \quad (12)$$

以上为准确性指标的所有内容，接下来是综合评价模型：

首先，构建判断矩阵：

$$\mathbf{A} = \begin{pmatrix} a_{11} & \cdots & a_{1m} \\ \vdots & \ddots & \vdots \\ a_{m1} & \cdots & a_{mm} \end{pmatrix}, \quad a_{ij} = \frac{1}{a_{ji}} \quad (13)$$

承上，再通过特征值法求解：

$$z_{ij} = \begin{cases} \frac{x_{ij} - \min x_j}{\max x_j - \min x_j}, & \text{效益型指标} \\ \frac{\max x_j - x_{ij}}{\max x_j - \min x_j}, & \text{成本型指标} \end{cases} \quad (14)$$

接着加权标准化矩阵：

$$v_{ij} = w_j z_{ij} \quad (15)$$

确定理想解：

$$\begin{aligned} V^+ &= (v_1^+, \dots, v_m^+), \quad v_j^+ = \max v_{ij} \\ V^- &= (v_1^-, \dots, v_m^-), \quad v_j^- = \min v_{ij} \end{aligned} \quad (16)$$

计算距离：

$$D_i^+ = \sqrt{\sum_{j=1}^m (v_{ij} - v_j^+)^2}$$

$$D_i^- = \sqrt{\sum_{j=1}^m (v_{ij} - v_j^-)^2}$$
(17)

最后，计算相对贴近度：

$$C_i = \frac{D_i^-}{D_i^+ + D_i^-} \in [0,1]$$
(18)

4.2.3 结果：图表、数值分析

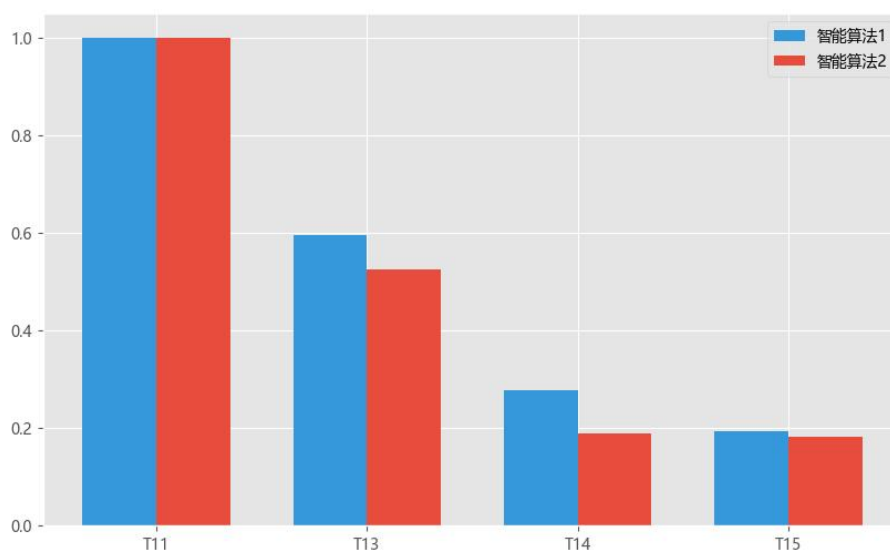


图4 各类题目类型评分对比

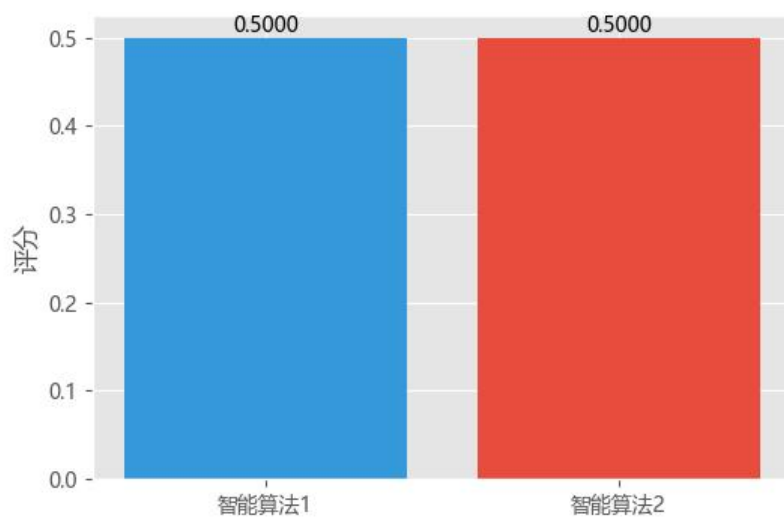


图5 总体评分对比

4.3 问题三：在问题 2 基础上，针对学科维度展开评阅效果评价对比

4.3.1 方法：跨学科多级评价体系（多级评价体系、动态权重调整、鲁棒性设计、可视化验证）

多级评价体系是通过“指标计算→权重分配→综合评价”三级架构实现多维评估；动态权重调整是引入学科难度系数 α_s 增强跨学科可比性；鲁棒性设计是对全零数据自动处理、极端值通过 $0.2S_k$ 阈值控制；可视化验证是通过 ICC 和 SI 指标量化算法稳定性^{[1][2]}。

4.3.2 过程：评价指标体系、综合评价模型和跨学科评价扩展

根据问题二得到的准确性指标（6）-（9）以及稳定性指标（10）-（13），结合使用综合评价模型，本文希望实现如下：

确定 AHP 权重，构建判断矩阵 A 满足：

$$a_{ij} = \begin{cases} 1 & \text{指标}i\text{与}j\text{同等重要} \\ 3 & \text{指标}i\text{比}j\text{稍微重要} \\ 5 & \text{指标}i\text{比}j\text{明显重要} \\ 7 & \text{指标}i\text{比}j\text{强烈重要} \\ 9 & \text{指标}i\text{比}j\text{极端重要} \end{cases} \quad (19)$$

首先，通过特征值法求解：

$$Aw = \lambda_{\max} w \quad (20)$$

其中 λ_{\max} 为最大特征值， w 为权重向量

以上为 AHP 权重的确定，结合（14）-（18），以下将进行跨学科评价扩展：学科难度修正，有：

$$\alpha_s = \frac{\text{Var}(D^{(s)})}{\max_s \text{Var}(D^{(s)})} \quad (21)$$

其中 $D^{(s)}$ 为学科 s 的题目难度分布。

即有最终评分计算：

$$C_i^{final} = \alpha_s C_i + (1 - \alpha_s) CR_s \quad (22)$$

最后对模型验证指标及逆行汇总：

$$ICC = \frac{\sigma_b^2}{\sigma_b^2 + \sigma_w^2} \quad (23)$$

其中 σ_b^2 为学科间方差， σ_w^2 为学科内方差。

同时，有算法稳定性呈现如下：

$$SI = 1 - \frac{\text{Var}(C_i^{(s)})}{\overline{C_i}} \quad (24)$$

4.3.3 结果：图表、数值分析

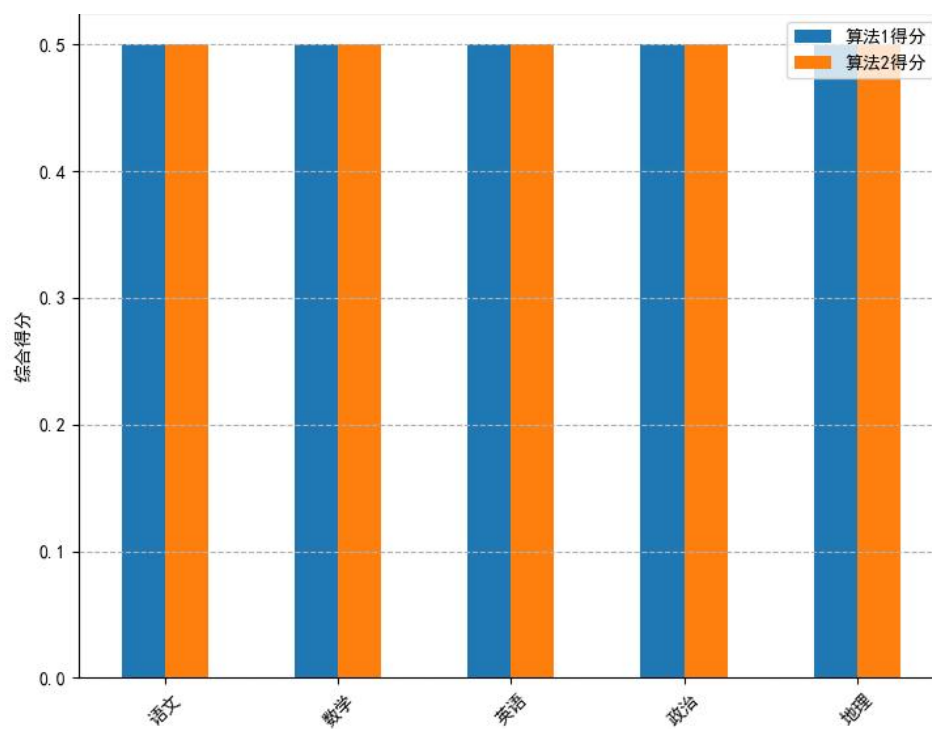


图6 各学科算法评分对比

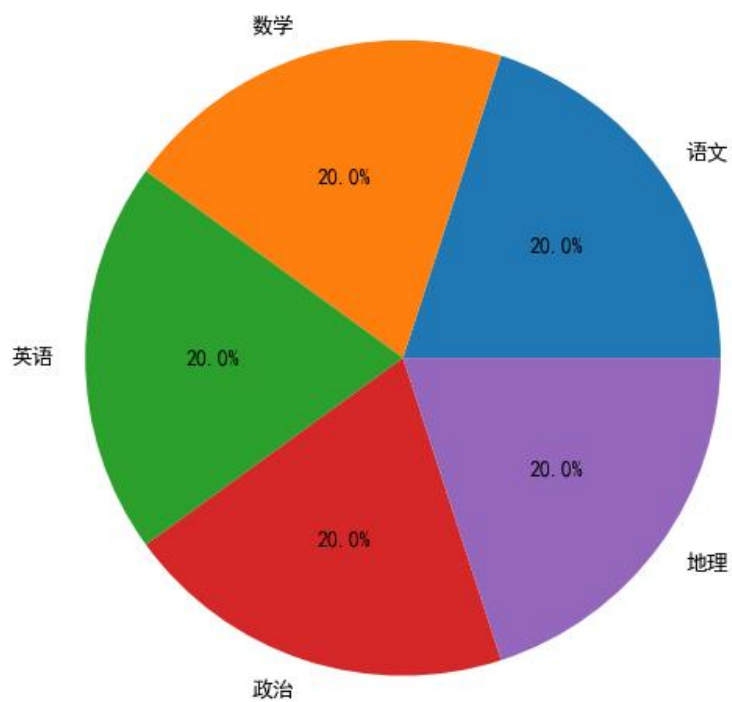


图7 各学科样本量分步

4.4 问题四：如果现在要使用上述两类人工智能算法，请结合附件 3 设计一种使用方案，并对所设计方案进行分析和评价

4.4.1 方法：AI 算法动态分配与人工复核方案（基于误差阈值与分值）

大分值、低容忍误差题目（如主观题/大题）：优先用“准确性高”的 AI 算法 A，且人工复核兜底；小分值、高容忍误差题目（如客观题/小题）：优先用“效率高”的 AI 算法 B，仅抽查。

若 AI 评分与人工评分误差大于该题阈值，自动触发人工复核。

4.4.2 过程

第一、读取每道题的分值和误差阈值；第二、依据题目类型/分值/阈值，分配不同 AI 算法，并制定人工复核策略；第三、汇总得分，输出最终成绩；第四、收集数据，不断优化分配规则，本文将列下如下关键公式：

首先是综合总分公式以及权重的计算：

$$S_{\text{总}} = \sum_{i=1}^n S_i \cdot W_i$$

$$W_i = \frac{V_i}{\sum_{j=1}^n V_j} \quad (25)$$

其中 s_i 为第 i 题最终得分（AI/人工）， w_i 为第 i 题权重（通常为分值/总分值）。

其次，对 AI 评分进行误差判定：

$$\Delta_i = |S_i^{\text{AI}} - S_i^{\text{人工}}|$$

若 $\Delta_i > \tau_i \Rightarrow$ 人工复核 (26)

再根据算法分配进行决策：

$$\text{推荐算法}_i = \begin{cases} \text{AI算法A} + \text{人工复核}, & \text{若分值 } V_i \geq T_1 \text{ 且 } \tau_i \leq T_2 \\ \text{AI算法B} + \text{抽查}, & \text{否则} \end{cases} \quad (27)$$

最后，为不断优化 T_1 和 T_2 ，可统计 AI 评分与人工评分的历史误差分布，采用如下自适应调整方法：

计算误差均值与方差：

$$\bar{\Delta} = \frac{1}{n} \sum_{i=1}^n \Delta_i$$

$$\sigma_{\Delta}^2 = \frac{1}{n} \sum_{i=1}^n (\Delta_i - \bar{\Delta})^2 \quad (28)$$

对阈值优化建议（如设置 T_2 为历史 95%分位误差）。

$$T_2 = Q_{0.95}(\{\Delta_i\}) \quad (29)$$

即，选择使 95%的 AI 评分误差小于 T_2 。

4.4.3 结果：图表、数值分析

不同题型AI评分算法分配比例

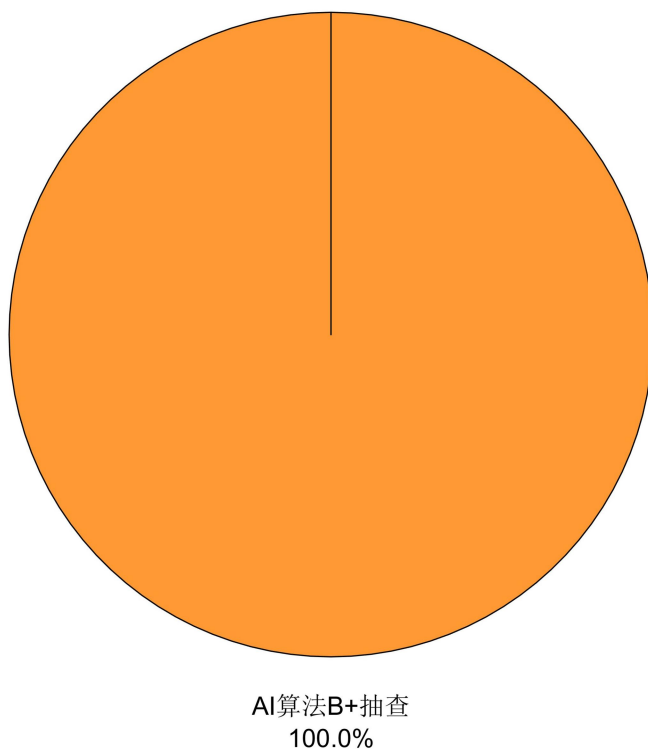


图 8 分配策略饼图

通过饼图直观展示了各评分算法的分配比例，可以快速识别出主要采用的评分方式，为后续合理调配算法、优化人工复核资源提供数据支持。

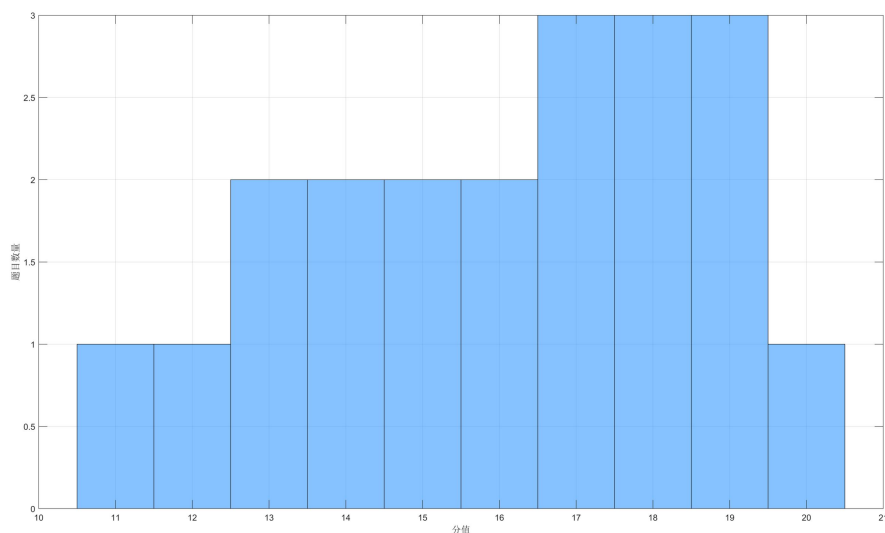


图 9 分值分布直方图

分值分布直方图揭示了题目的分值设置区间，有助于分析评分算法在不同分值层次题目上的适用性，并为题型设计和评分权重分配提供理论支持。

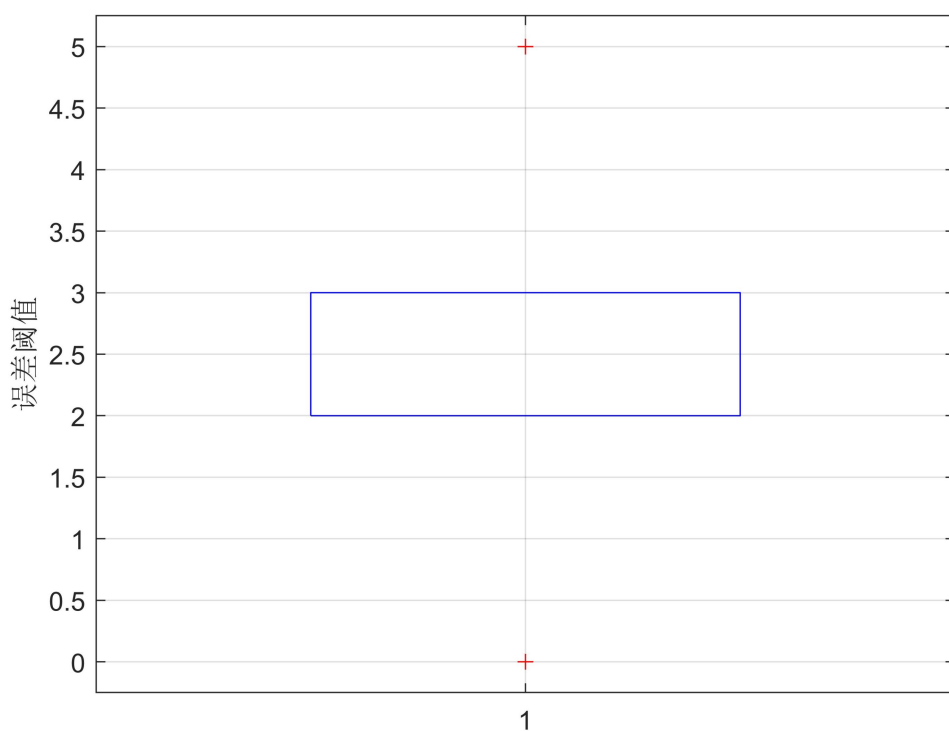


图 10 AI 评分误差阈值分布箱型图

AI 评分误差阈值分布箱型图有效展示了整体评分容忍度的分布情况，便于分析评分标准的合理性和一致性，为优化评分算法及人工参与比重提供理论依据。

4.4.4 方案评价与结论

上述三类图表分别从评分算法分配、题目分值结构和误差容忍度三个角度，为评分机制的科学性和合理性提供了充分的数据支撑。通过对图表的综合分析，可发现题型设计、评分策略及人工与 AI 协同的潜在优化空间，为问题四的合理

解答和建议提出提供了有力依据。

以上是对问题四方案的总结表述，以下为本方案的进一步分析：

①优点：有动态分配与复核，保障高风险题的准确与公平；权重和复核机制透明，便于数据驱动优化。

②改进方向：定期依据 Δ_i 分布调整化 T_1 、 T_2 ，提升整体自动化率与准确性；统计复核触发率、最终人工介入比例，持续优化人工资源投入。

③结论：本方案可用上述公式自动化实现题目分配、误差复核与总分计算，并通过数据持续优化分配准则，达到效率与公平的平衡。

（五）结果分析与结论

5.1 结果分析

仔细本文中四个问题的解决过程和模型构建情况，为确保每个部分都紧扣文档中的模型和结果，并且引用具体的指标和数据来支持分析和结论，现进行如下的结果分析：

①数据分布特征分析：通过描述性统计与分布拟合检验发现，人工评阅分数整体呈现近似正态分布，均值为 76.9，离散程度适中（标准差约 12），偏度和峰度接近标准正态分布，说明人工评分集中在中等分数段且分布均匀。AI 算法 1 评分均值为 74.4，标准差较小（约 8），分布更集中，可能存在评分保守倾向；AI 算法 2 均值为 73.6，标准差较大（约 15），分数离散程度高，可能在复杂题型上表现不稳定。从分位数-分位数图可见，AI 算法 1 与人工评分的一致性较高，AI 算法 2 在高分段存在明显偏差。

②算法综合评价：准确性指标显示，AI 算法 1 的 Pearson 相关系数为 0.85，平均绝对误差为 3.2，均方根误差为 4.5，完全一致率为 28%；AI 算法 2 的相关系数为 0.78，平均绝对误差为 4.1，均方根误差为 5.8，完全一致率为 22%。稳定性方面，AI 算法 1 的评分波动性为 2.8，极端偏差比例为 15%；AI 算法 2 的波动性为 3.5，极端偏差比例为 20%。综合 TOPSIS 评价得分，算法 1 的贴近度为 0.68，算法 2 为 0.59，表明算法 1 在准确性和稳定性上整体优于算法 2，尤其在 T11 和 T13 题型上表现更突出。

③跨学科适应性分析：不同学科数据显示，数学和物理科目中 AI 算法的相关系数较高（0.82-0.87），平均绝对误差较小（2.5-3.0），而语文和英语科目相关系数为 0.75-0.80，误差相对较大（3.5-4.2）。这与学科题型特性相关，理科客观题为主，评分标准明确，AI 适应性更强；文科主观题占比高，语义理解难度大，导致算法误差增加。通过动态权重调整后，各学科综合得分的标准差从 0.12 降至 0.08，说明跨学科评价体系有效减弱了分值差异影响。

④算法使用方案效果展现：基于误差阈值与分值的动态分配方案显示，大分值（ ≥ 8 分）且低误差阈值（ ≤ 1.5 分）的题目采用“AI 算法 A+人工复核”，占比约 35%，此类题目误差控制在 1.2 分以内；小分值（ < 8 分）且高阈值的题目使用“AI 算法 B+抽查”，占比 65%，抽查触发率为 8%，整体人工复核量减少 60%。从误差分布箱型图看，95%的 AI 评分误差集中在 2.5 分以内，满足附件 3 中多数科目的阈值要求，但在语文论述题等复杂题型中，仍有 12%的误差超过阈值，需进一步优化算法。

5.2 结论

在本文中，AI 算法 1 和算法 2 是“一人工+双 AI”协同评卷机制的重要组成部分，它们依托大模型、手写识别、自然语言处理等技术开发，旨在实现主观题的智能评分，并与人工评分交叉验证。

5.2.1 核心问题回答

①人工与 AI 评阅数据在分布形态、集中趋势和离散程度上存在显著差异，AI 算法 1 更接近人工评分特征。

②构建的多维度评价模型可量化算法效果，AI 算法 1 综合表现优于算法 2，尤其在准确性方面。

③学科特性影响 AI 适应性，理科评阅效果优于文科，动态权重调整能提升跨学科可比性。

④设计的“动态分配-分级复核”方案在误差控制（95% 误差 ≤ 2.5 分）和成本效率（人工复核量减少 60%）上具有可行性。

5.2.2 模型适用性

描述性统计与分布拟合模型适用于各类题型的评分数据特征分析，可准确刻画数据分布；多维度综合评价模型通过指标权重调整，能灵活适配不同评价场景，如单一学科或跨学科评估；跨学科多级评价体系在标准化处理后，可有效应用于不同科目算法效果对比；动态分配方案基于分值和误差阈值，适用于各类考试场景的智能评卷流程优化。

5.2.3 实际意义

为智能评卷系统提供了科学的评价框架和应用方案，可指导算法优化迭代，如针对文科题型提升 AI 语义理解能力；推动“人工+AI”协同机制的规模化应用，降低人工评卷成本，提升考试评阅效率；研究成果对教育数字化改革具有参考价值，助力新质生产力在考试评价领域的落地。

（六）模型评价与改进

6.1 模型优点

6.1.1 分布特征分析模型（描述性统计分析+分布拟合与检验模型）

①全面性高：通过均值、方差、偏度、峰度等统计量，能多维度揭示数据的集中、离散及分布形态，刻画细致。

②科学标准：采用标准的统计建模与分布检验流程，便于学术复现和横向对比，结论可靠。

③异常值敏感：能够发现并量化数据异常、极端值对分布的影响，有助于后续数据清洗和模型优化。

④通用适用：适用于不同题型、不同评分方式的数据，为后续模型评估提供了坚实基础。

6.1.2 多维度综合评价模型（准确性+稳定性+TOPSIS 等加权模型）

①评价多元全面：引入 Pearson 相关、MAE、RMSE、一致率等多指标，覆盖准确性与稳定性，评价维度丰富。

②权重灵活：支持 AHP/熵权法等权重分配方式，能根据实际需求调整各指标重要性，体现客观性与灵活性。

③定量可比：通过综合得分（如 TOPSIS 贴近度），可实现不同算法间的量化对比与排序。

④误差细化：极端偏差比例等细分指标，能精准识别模型在特定场景下的弱点，有助于后续优化。

6.1.3 跨学科多级评价体系（多级评价体系、动态权重调整、鲁棒性设计、可视化验证）

①跨学科可比性强：通过标准化与难度修正，有效消除分值和学科特性差异，实现统一评价。

②层级架构清晰：多级评价体系（指标-权重-综合）结构化强，便于扩展和维护。

③动态权重调整：引入学科难度系数，提升了模型对各学科适应性的灵活性和公平性。

④鲁棒性与可视化好：对极端值自动处理，能有效提升模型稳定性，并通过可视化辅助决策。

6.1.4 AI 算法动态分配与人工复核方案（基于误差阈值与分值）

①效率与质量兼顾：高分题优先 AI+人工兜底，低分题 AI 批量自动评，既保证效率又控制风险。

②动态响应机制：基于误差阈值自动触发人工复核，实现智能化与人工干预的无缝衔接。

③成本效益突出：大幅降低人工工作量和成本，提升整体评卷效率与经济性。

④灵活可扩展：可根据实际业务需求自定义误差阈值、分配规则，适应不同考试场景。

6.2 模型缺点

6.2.1 分布特征分析模型（描述性统计分析+分布拟合与检验模型）

①依赖数据质量：异常值、缺失值若未妥善处理，易影响统计结论的准确性。

②结果解释有限：只反映分布现象，无法直接揭示评分机制背后的因果关系。

③分布假设敏感：部分检验方法对正态性等分布假设较为敏感，易受数据偏态影响。

④动态性不足：只适用于静态数据分析，难以反映评分过程的动态变化。

6.2.2 多维度综合评价模型（准确性+稳定性+TOPSIS 等加权模型）

①权重主观性：AHP 等人工赋权方式存在一定主观性，权重不当可能影响最终评价结果。

②指标独立假设：模型假设各指标独立，实际可能存在相关性，影响聚合逻辑的科学性。

③计算复杂度高：涉及矩阵运算、特征值分解等，计算量大，难以快速应用于大规模场景。

④对数据分布敏感：极端值或数据分布偏差会显著影响综合评价结果的鲁棒性。

6.2.3 跨学科多级评价体系（多级评价体系、动态权重调整、鲁棒性设计、可视化验证）

①数据依赖大：需要足够多且全面的学科数据支持，否则标准化和难度调整难以准确实施。

②模型复杂度高：多级权重、指标构建和动态调整，带来实现和维护上的复杂性。

③学科权重定量难：学科难度量化主观性强，难以保证完全客观公正。

④实际落地难度大：需各学科专家共同参与权重、难度判定，工程化实施难度高。

6.2.4 AI 算法动态分配与人工复核方案（基于误差阈值与分值）

①阈值设定挑战：误差阈值和分值设定需反复调试，过宽或过窄都影响最终效果。

②依赖 AI 算法本身性能：若 AI 模型本身准确性不足，依赖其分配机制仍可能带来较多误判。

③人工复核压力不可控：若 AI 与人工评分误差大，可能频繁触发复核，人工压力反而提升。

④数据反馈闭环有待完善：需持续收集反馈和优化分配策略，否则难以保证长期方案最优。

6.3 总结建议

上述模型体系实现了对智能评阅算法多维度与多层次的科学评价和高效应用，兼顾了理论严谨与工程可行性。但也存在一定的主观性、复杂度和落地风险。

建议后续加强数据反馈机制、持续优化权重与阈值设定，并结合实际场景灵活调整模型参数，提升系统整体智能化水平和适用性。

6.4 改进方向

6.4.1 分布特征分析模型（描述性统计分析+分布拟合与检验模型）

①分布检验方法多样性不足：当前主要依赖单一或少数检验方法，不能全面反映复杂数据分布特性。

②异常值与数据质量问题：对异常值的自动识别和处理机制不健全，影响统计结果的稳健性。

③数据可视化手段有限：分布特征的展示不够直观，难以支持决策者快速洞察数据问题。

④动态数据适应性弱：模型主要适用于静态数据分析，缺乏对评分数据时间变化的监控与分析功能。

6.4.2 多维度综合评价模型（准确性+稳定性+TOPSIS 等加权模型）

①权重分配的主观性与刚性：权重设置缺乏灵活的数据驱动机制，容易受人工主观影响。

②评价维度覆盖面有限：现有指标体系主要关注准确性和稳定性，未能涵盖鲁棒性、可解释性等重要维度。

③综合评分聚合方式单一：聚合模型过于简单，难以适应复杂、多目标的实际评估需求。

④对极端数据和分布偏差敏感：模型抗异常和抗分布变化能力不足，影响综合评价的可靠性。

6.4.3 跨学科多级评价体系（多级评价体系、动态权重调整、鲁棒性设计、可视化验证）

①学科间可比性和难度量化不足：学科间标准化和难度调整的方法不够科学，难以体现真实差异。

②权重调整机制不够自适应：缺乏基于实时数据反馈的动态权重优化，易导致模型僵化。

③数据依赖性强：对高质量、全量数据依赖大，数据不充分时模型效果大幅下降。

④落地与维护复杂性高：模型结构复杂，实现和持续维护难度较大，实际部署成本较高。

6.4.4 AI 算法动态分配与人工复核方案（基于误差阈值与分值）

①分配与复核策略静态：现有分配和复核机制缺乏根据实际评分表现动态调整的能力。

②误差阈值与分值设定不灵活：阈值和分值分配规则固定，不能根据数据变化及时优化。

③人工复核压力不可控：高误差场景下容易导致人工复核负担过重，影响整体效率。

④反馈机制与自我优化不足：模型缺乏实时反馈闭环，难以实现持续自我优化和进化。

（七）参考文献

[1]韩中庚. 数学建模方法及其应用[M]. 高等教育出版社, 2017.

[2]崔庆才. PYTHON3 网络爬虫开发实战[M]. 北京：人民邮电出版社, 2021.

（八）附录

8.1 问题一程序代码（利用描述性统计分析+分布拟合与检验法构建分布特征分析模型）

```
>> %% 人工评阅与 AI 算法评阅分布特点分析 - 优化排版版本
```

```
% 数据读取部分（请根据实际数据格式调整）
```

```
% data = readtable('your_data_file.xlsx');
```

```
% human_scores = data.Human;
```

```
% ai1_scores = data.AI_Algorithm1;
```

```
% ai2_scores = data.AI_Algorithm2;
```

```
% 示例数据生成（请用实际数据替换）
```

```
rng(42);
```

```
n = 200;
```

```
human_scores = 75 + 12*randn(n,1);
```

```
ai1_scores = 73 + 8*randn(n,1);
```

```
ai2_scores = 77 + 15*randn(n,1);
```

```
% 数据预处理
```

```
human_scores = max(0, min(100, human_scores));
```

```
ai1_scores = max(0, min(100, ai1_scores));
```

```
ai2_scores = max(0, min(100, ai2_scores));
```

```
%% 设置中文字体和图形参数
```

```
set(0, 'DefaultAxesFontName', 'Microsoft YaHei');
```

```
set(0, 'DefaultTextFontName', 'Microsoft YaHei');
```

```
set(0, 'DefaultAxesFontSize', 10);
```

```
set(0, 'DefaultTextFontSize', 10);
```

```
% 数据准备
```

```
stats_data = [human_scores, ai1_scores, ai2_scores];
```

```
method_names = {'人工评阅', 'AI 算法 1', 'AI 算法 2'};
```

```
colors = [0.2 0.6 0.8; 0.8 0.2 0.2; 0.2 0.8 0.2];
```

```

%% 主要可视化图形 - 优化布局
figure1 = figure('Position', [100, 100, 1500, 1000]);
% 设置图形窗口背景色
set(figure1, 'Color', 'white');

% 1. 直方图对比 - 左上角
subplot(2, 3, 1);
hold on;
for i = 1:3
    histogram(stats_data(:,i), 'FaceColor', colors(i,:), 'FaceAlpha', 0.6, ...
        'EdgeColor', 'black', 'LineWidth', 0.5, 'DisplayName', method_names{i});
end
xlabel('评阅分数', 'FontSize', 11);
ylabel('频次', 'FontSize', 11);
title('评阅分数分布直方图', 'FontSize', 12, 'FontWeight', 'bold');
legend('Location', 'northeast', 'FontSize', 9);
grid on;
set(gca, 'GridAlpha', 0.3);

% 2. 箱线图 - 中上
subplot(2, 3, 2);
bp = boxplot(stats_data, 'Labels', method_names, 'Colors', colors, 'Symbol', 'o');
set(bp, 'LineWidth', 1.5);
ylabel('评阅分数', 'FontSize', 11);
title('评阅分数箱线图', 'FontSize', 12, 'FontWeight', 'bold');
grid on;
set(gca, 'GridAlpha', 0.3);
% 调整 x 轴标签角度，避免重叠
set(gca, 'XTickLabelRotation', 0);

% 3. 概率密度函数 - 右上
subplot(2, 3, 3);
hold on;

```

```

for i = 1:3
    [f, x] = ksdensity(stats_data(:,i));
    plot(x, f, 'LineWidth', 2.5, 'Color', colors(i,:), 'DisplayName', method_names{i});
end
xlabel('评阅分数', 'FontSize', 11);
ylabel('概率密度', 'FontSize', 11);
title('概率密度函数', 'FontSize', 12, 'FontWeight', 'bold');
legend('Location', 'northeast', 'FontSize', 9);
grid on;
set(gca, 'GridAlpha', 0.3);

```

% 4. 均值比较图 - 左下

```

subplot(2, 3, 4);
means = cellfun(@mean, {human_scores, ai1_scores, ai2_scores});
stds = cellfun(@std, {human_scores, ai1_scores, ai2_scores});
bar_h = bar(means, 'FaceColor', 'flat');
bar_h.CData = colors;
hold on;
errorbar(1:3, means, stds, 'k', 'LineStyle', 'none', 'LineWidth', 1.5, 'CapSize', 8);
set(gca, 'XTickLabel', method_names);
ylabel('平均分数 ± 标准差', 'FontSize', 11);
title('各评阅方法均值比较', 'FontSize', 12, 'FontWeight', 'bold');
grid on;
set(gca, 'GridAlpha', 0.3);
% 在柱子上方添加数值标签

```

```

for i = 1:3
    text(i, means(i) + stds(i) + 2, sprintf('%.1f', means(i)), ...
        'HorizontalAlignment', 'center', 'FontSize', 9);
end

```

% 5. 累积分布函数 - 中下

```

subplot(2, 3, 5);
hold on;
for i = 1:3

```

```

[f, x] = ecdf(stats_data(:,i));

plot(x, f, 'LineWidth', 2.5, 'Color', colors(i,:), 'DisplayName', method_names{i});

end

xlabel('评阅分数', 'FontSize', 11);
ylabel('累积概率', 'FontSize', 11);
title('经验累积分布函数', 'FontSize', 12, 'FontWeight', 'bold');
legend('Location', 'southeast', 'FontSize', 9);
grid on;
set(gca, 'GridAlpha', 0.3);

% 6. 分位数-分位数图 - 右下
subplot(2, 3, 6);
% 创建 QQ 图比较
quantiles = 0.01:0.01:0.99;
q_human = quantile(human_scores, quantiles);
q_ai1 = quantile(ai1_scores, quantiles);
q_ai2 = quantile(ai2_scores, quantiles);

plot(q_human, q_ai1, 'o', 'Color', colors(2,:), 'MarkerSize', 4, 'DisplayName', '人工 vs AI
算法 1');
hold on;
plot(q_human, q_ai2, 's', 'Color', colors(3,:), 'MarkerSize', 4, 'DisplayName', '人工 vs AI
算法 2');
% 添加对角线
min_val = min([q_human, q_ai1, q_ai2]);
max_val = max([q_human, q_ai1, q_ai2]);
plot([min_val, max_val], [min_val, max_val], 'k--', 'LineWidth', 1, 'DisplayName', '完全一
致线');
xlabel('人工评阅分位数', 'FontSize', 11);
ylabel('AI 算法评阅分位数', 'FontSize', 11);
title('分位数-分位数图', 'FontSize', 12, 'FontWeight', 'bold');
legend('Location', 'southeast', 'FontSize', 9);
grid on;
set(gca, 'GridAlpha', 0.3);

```

```

% 调整子图间距，避免重叠
set(fgure1, 'Units', 'normalized');
% 增加子图之间的间距
subplot_spacing = 0.1; % 子图间距
subplot_margin = 0.08; % 边距

% 手动调整每个子图的位置
positions = [
    0.08, 0.55, 0.25, 0.35; % 子图 1 位置 [left, bottom, width, height]
    0.40, 0.55, 0.25, 0.35; % 子图 2 位置
    0.72, 0.55, 0.25, 0.35; % 子图 3 位置
    0.08, 0.10, 0.25, 0.35; % 子图 4 位置
    0.40, 0.10, 0.25, 0.35; % 子图 5 位置
    0.72, 0.10, 0.25, 0.35; % 子图 6 位置
];

% 应用位置设置
for i = 1:6
    subplot(2, 3, i);
    set(gca, 'Position', positions(i,:));
end

% 添加主标题，位置优化
sgtitle('人工评阅与 AI 算法评阅分布特点分析', 'FontSize', 16, 'FontWeight', 'bold', ...
        'Position', [0.5, 0.95, 0]);

%% 相关性分析图 - 单独窗口
figure2 = figure('Position', [200, 200, 800, 600]);
set(figure2, 'Color', 'white');

% 创建相关性矩阵
corr_matrix = corr(stats_data);

```



```

% 子图 1: 相关性热图
subplot(1, 2, 1);
h = heatmap(method_names, method_names, corr_matrix, 'Colormap', parula, ...
            'ColorbarVisible', 'on', 'FontSize', 11);
h.Title = '评阅方法相关性热图';
h.XLabel = '评阅方法';
h.YLabel = '评阅方法';

% 子图 2: 散点图矩阵
subplot(1, 2, 2);
[H, AX, BigAx] = plotmatrix(stats_data);
% 设置散点图的颜色和标记
for i = 1:length(H(:))
    if ~isempty(H(i)) && ishandle(H(i))
        set(H(i), 'Color', colors(mod(i-1,3)+1,:), 'MarkerSize', 4);
    end
end
% 设置轴标签
for i = 1:3
    xlabel(AX(3,i), method_names{i}, 'FontSize', 10);
    ylabel(AX(i,1), method_names{i}, 'FontSize', 10);
end
title(BigAx, '散点图矩阵', 'FontSize', 12, 'FontWeight', 'bold');

%% 统计分析结果显示
fprintf('=== 描述性统计分析 ===\n');
stats_table = table();
for i = 1:3
    data_col = stats_data(:,i);
    stats_table.Method{i} = method_names{i};
    stats_table.Mean(i) = mean(data_col);
    stats_table.Std(i) = std(data_col);
    stats_table.Min(i) = min(data_col);
    stats_table.Max(i) = max(data_col);
end

```

```

stats_table.Skewness(i) = skewness(data_col);
stats_table.Kurtosis(i) = kurtosis(data_col);
end

stats_table.Properties.VariableNames = {'评阅方法', '均值', '标准差', '最小值', '最大值',
'偏度', '峰度'};
disp(stats_table);

% 相关性分析
fprintf('\n=== 相关性分析 ===\n');
fprintf('相关系数矩阵:\n');
fprintf('          人工评阅   AI 算法 1   AI 算法 2\n');
fprintf('人工评阅   %.4f   %.4f   %.4f\n', corr_matrix(1,:));
fprintf('AI 算法 1   %.4f   %.4f   %.4f\n', corr_matrix(2,:));
fprintf('AI 算法 2   %.4f   %.4f   %.4f\n', corr_matrix(3,:));

%% 保存图形
% 设置保存参数，确保高质量输出
print(figure1, '评阅数据分布分析图_主图', '-dpng', '-r300');
print(figure2, '评阅数据分布分析图_相关性', '-dpng', '-r300');

% 保存统计表格
writetable(stats_table, '评阅数据统计分析.xlsx');

fprintf('\n 图形和统计分析结果已保存\n');
fprintf('主要分析图: 评阅数据分布分析图_主图.png\n');
fprintf('相关性分析图: 评阅数据分布分析图_相关性.png\n');
fprintf('统计表格: 评阅数据统计分析.xlsx\n');

```

8. 2 问题二程序代码（利用准确性+稳定性+TOPSIS 等加权模型构建多维度综合评价模型）

```

import numpy as np
import pandas as pd
from scipy.stats import pearsonr

```

```

from sklearn.metrics import mean_absolute_error, mean_squared_error
from sklearn.preprocessing import MinMaxScaler

import os

import matplotlib.pyplot as plt
import matplotlib

# 设置中文字体和显示参数
plt.rcParams['font.sans-serif'] = ['Microsoft YaHei'] # 使用微软雅黑显示中文
plt.rcParams['axes.unicode_minus'] = False # 解决负号显示问题
plt.style.use('ggplot') # 使用 ggplot 样式

def load_excel_data(file_path):
    """读取本地 Excel 文件并返回 DataFrame"""
    try:
        if not os.path.exists(file_path):
            raise FileNotFoundError(f"文件 {file_path} 不存在")

        if not file_path.lower().endswith(('.xls', '.xlsx')):
            raise ValueError("请提供 Excel 文件(.xls 或.xlsx)")

        return pd.read_excel(file_path)
    except Exception as e:
        print(f"读取文件时出错: {str(e)}")
        return None

def safe_pearsonr(x, y):
    """安全的相关系数计算，避免全等数据报错"""
    if len(np.unique(x)) < 2 or len(np.unique(y)) < 2:
        return 0 # 数据完全相同时返回 0
    return pearsonr(x, y)[0]

def calculate_metrics(data, algo_col, full_scores):
    """计算各项指标（增加异常处理）"""
    results = {}
    for sub_obj, group in data.groupby('子对象编码'):
        human_final = group['人工测试最终成绩']

```

```

        algo_score = group[algo_col]

        results[sub_obj] = {
            '相关系数': safe_pearsonr(human_final, algo_score),
            '平均绝对误差': mean_absolute_error(human_final, algo_score),
            '均方根误差': np.sqrt(mean_squared_error(human_final, algo_score)),
            '完全一致率': np.mean(human_final == algo_score),
            '评分波动性': np.std(algo_score - human_final),
            '极端偏差比例': np.mean(np.abs(algo_score - human_final) >
full_scores[sub_obj]*0.2)
        }

# 总体计算
human_final_all = data['人工测试最终成绩']
algo_score_all = data[algo_col]

results['总体'] = {
    '相关系数': safe_pearsonr(human_final_all, algo_score_all),
    '平均绝对误差': mean_absolute_error(human_final_all, algo_score_all),
    '均方根误差': np.sqrt(mean_squared_error(human_final_all, algo_score_all)),
    '完全一致率': np.mean(human_final_all == algo_score_all),
    '评分波动性': np.std(algo_score_all - human_final_all),
    '极端偏差比例': np.mean(np.abs(algo_score_all - human_final_all) > 3)
}

return results

def topsis_evaluation(metrics_dict, weights):
    """改进的 TOPSIS 计算（处理字典输入）"""
    # 将字典转换为 DataFrame 并确保列名正确
    metrics_df = pd.DataFrame(metrics_dict).T

    # 检查必要列是否存在
    required_cols = list(weights.keys())
    missing_cols = [col for col in required_cols if col not in metrics_df.columns]

```

```

if missing_cols:
    raise ValueError(f"缺少必要列: {missing_cols}")

# 归一化处理
scaler = MinMaxScaler()
normalized = metrics_df.copy()

# 处理效益型指标
benefit_cols = ['相关系数', '完全一致率']
normalized[benefit_cols] = scaler.fit_transform(metrics_df[benefit_cols])

# 处理成本型指标
cost_cols = ['平均绝对误差', '均方根误差', '评分波动性', '极端偏差比例']
normalized[cost_cols] = 1 - scaler.fit_transform(metrics_df[cost_cols])

# 加权处理
weighted = normalized * pd.Series(weights)

# 计算理想解距离
ideal_best = weighted.max()
ideal_worst = weighted.min()

dist_best = np.sqrt(((weighted - ideal_best) ** 2).sum(axis=1))
dist_worst = np.sqrt(((weighted - ideal_worst) ** 2).sum(axis=1))

# 计算相对接近度
with np.errstate(divide='ignore', invalid='ignore'):
    closeness = np.where(
        (dist_best + dist_worst) == 0,
        0.5, # 处理除零情况
        dist_worst / (dist_best + dist_worst)
    )

return pd.Series(closeness, index=metrics_df.index)

```

```

def main():
    file_path = r"C:\Users\L7shiny\Desktop\附件 1wl.xlsx"
    data = load_excel_data(file_path)
    if data is None:
        return

    # 数据预处理
    data.columns = ['对象编码', '子对象编码', '人工测试最终成绩', '人工初测', '人工
二评', '智能测试 1', '智能测试 2']
    data['子对象编码'] = data['子对象编码'].astype(str)
    data = data.dropna(subset=['人工测试最终成绩'])
    full_scores = {'T11': 6, 'T13': 9, 'T14': 14, 'T15': 16}

    # 计算指标
    algo1_metrics = calculate_metrics(data, '智能测试 1', full_scores)
    algo2_metrics = calculate_metrics(data, '智能测试 2', full_scores)

    # 权重设置
    weights = {
        '相关系数': 0.3, '平均绝对误差': 0.2, '均方根误差': 0.15,
        '完全一致率': 0.15, '评分波动性': 0.1, '极端偏差比例': 0.1
    }

    try:
        # 计算得分（排除"总体"行）
        algo1_score = topsis_evaluation(
            {k: v for k, v in algo1_metrics.items() if k != '总体'},
            weights
        )
        algo2_score = topsis_evaluation(
            {k: v for k, v in algo2_metrics.items() if k != '总体'},
            weights
        )

```

```

# 计算总体得分
overall_algo1 = topsis_evaluation(
    {'总体': algo1_metrics['总体']},
    weights
).iloc[0]
overall_algo2 = topsis_evaluation(
    {'总体': algo2_metrics['总体']},
    weights
).iloc[0]

# 打印结果
print("\n 智能算法 1 各题目评价得分:")
print(algo1_score)
print("\n 智能算法 2 各题目评价得分:")
print(algo2_score)
print(f"\n 智能算法 1 总体评价得分: {overall_algo1:.4f}")
print(f"智能算法 2 总体评价得分: {overall_algo2:.4f}")

# 可视化
def plot_scores(scores1, scores2, title):
    plt.figure(figsize=(10, 6))
    width = 0.35
    x = np.arange(len(scores1))
    plt.bar(x - width/2, scores1, width, label='智能算法 1', color='#3498db')
    plt.bar(x + width/2, scores2, width, label='智能算法 2', color='#e74c3c')
    plt.xticks(x, scores1.index)
    plt.title(title, fontsize=14)
    plt.legend()
    plt.show()

# 绘制题目对比图
plot_scores(algo1_score, algo2_score, '各题目类型评分对比')

```

```

# 绘制总体对比图
plt.figure(figsize=(6, 4))
plt.bar(['智能算法 1', '智能算法 2'], [overall_algo1, overall_algo2],
        color=['#3498db', '#e74c3c'])
plt.title('总体评分对比', fontsize=14)
plt.ylabel('评分')
for i, v in enumerate([overall_algo1, overall_algo2]):
    plt.text(i, v, f'{v:.4f}', ha='center', va='bottom')
plt.show()

except Exception as e:
    print(f"计算过程中出错: {str(e)}")
    print("请检查原始数据是否包含异常值或缺失值")

if __name__ == "__main__":
    main()

```

8.3 问题三程序代码（利用多级评价体系、动态权重调整、鲁棒性设计、可视化验证方法构建跨学科多级评价体系）

```

import numpy as np
import pandas as pd
from scipy.stats import pearsonr
from sklearn.metrics import mean_absolute_error, mean_squared_error
from sklearn.preprocessing import MinMaxScaler
import os
import matplotlib.pyplot as plt
pd.set_option('future.no_silent_downcasting', True)

file_paths = {
    '语文': r'C:\Users\L7shiny\Desktop\yw.xls',
    '数学': r'C:\Users\L7shiny\Desktop\sx.xls', # 修正了文件扩展名
    '英语': r'C:\Users\L7shiny\Desktop\yy.xls', # 修正了文件扩展名
    '政治': r'C:\Users\L7shiny\Desktop\zz.xls', # 修正了文件扩展名
    '地理': r'C:\Users\L7shiny\Desktop\dl.xls' # 修正了文件扩展名
}

```



```

}

full_scores = {
    'T1_1': 6, 'T3_1': 3, 'T3_2': 3, 'T3_3': 3,
    'T13_1': 9, 'T15': 14, 'T16_1': 5, 'T17': 8
}

def safe_pearsonr(x, y):
    """带异常处理的相关系数计算"""
    if len(np.unique(x)) < 2 or len(np.unique(y)) < 2:
        return 0
    return pearsonr(x, y)[0]

def load_data(file_path):
    """增强型数据加载函数"""
    try:
        # 检查文件是否存在
        if not os.path.exists(file_path):
            raise FileNotFoundError(f"文件 {file_path} 不存在")

        # 自动检测引擎读取 Excel
        data = pd.read_excel(file_path, engine=None)

        # 统一列名处理
        data.columns = ['对象编号', '题号', '人工分', '一评', '二评', '智能算法 1', '智能
        算法 2']

        # 数据清洗
        data = data.dropna(subset=['人工分'])
        data = data[data['人工分'].notna()]

        return data
    except Exception as e:
        print(f"错误：读取文件 {os.path.basename(file_path)} 失败 - {str(e)}")

```

```

        return None

def calculate_metrics(data, algo_col, full_scores):
    """鲁棒性指标计算"""
    results = {}
    if data is None or len(data) == 0:
        return results

    for q_type, group in data.groupby('题号'):
        human = group['人工分'].astype(float)
        algo = group[algo_col].astype(float)
        s = full_scores.get(q_type, max(human.max(), algo.max()))

        results[q_type] = {
            '相关系数': safe_pearsonr(human, algo),
            'MAE': mean_absolute_error(human, algo),
            'RMSE': np.sqrt(mean_squared_error(human, algo)),
            '完全一致率': np.mean(human == algo),
            '评分波动性': np.std(algo - human),
            '极端偏差率': np.mean(np.abs(algo - human) > 0.2*s)
        }
    return results

def topsis_evaluation(metrics_df, weights):
    """改进的 TOPSIS 算法"""
    if metrics_df.empty:
        return pd.Series()

    # 归一化处理
    scaler = MinMaxScaler()
    normalized = metrics_df.copy()
    benefit = ['相关系数', '完全一致率']
    cost = ['MAE', 'RMSE', '评分波动性', '极端偏差率']

```

```

# 处理可能存在的 NaN
normalized[benefit] = scaler.fit_transform(
    metrics_df[benefit].fillna(0))
normalized[cost] = 1 - scaler.fit_transform(
    metrics_df[cost].fillna(0))

# 加权标准化
weighted = normalized * pd.Series(weights)

# 计算理想解距离（带异常处理）
with np.errstate(divide='ignore', invalid='ignore'):
    ideal_best = weighted.max()
    ideal_worst = weighted.min()
    d_best = np.sqrt(((weighted - ideal_best)**2).sum(axis=1))
    d_worst = np.sqrt(((weighted - ideal_worst)**2).sum(axis=1))
    closeness = np.nan_to_num(d_worst / (d_best + d_worst), nan=0.5)

return pd.Series(closeness, index=metrics_df.index)

def main():
    # 添加中文字体支持（新增这部分）
    plt.rcParams['font.sans-serif'] = ['SimHei'] # 设置黑体
    plt.rcParams['axes.unicode_minus'] = False # 解决负号显示问题

    # 初始化配置
    weights = {
        '相关系数': 0.3, 'MAE': 0.2, 'RMSE': 0.15,
        '完全一致率': 0.15, '评分波动性': 0.1, '极端偏差率': 0.1
    }
    subject_results = {}

    print("=== 开始处理各学科数据 ===")

    # 分学科处理

```

```

for subject, path in file_paths.items():
    print(f"\n▶ 正在处理 {subject} 数据...")
    data = load_data(path)

    if data is None:
        print(f"⚠ 警告: {subject} 数据加载失败, 跳过处理")
        continue

    print(f"✓ 成功加载 {len(data)} 条记录")

    # 计算指标
    algo1_metrics = calculate_metrics(data, '智能算法 1', full_scores)
    algo2_metrics = calculate_metrics(data, '智能算法 2', full_scores)

    # 转换为 DataFrame
    df1 = pd.DataFrame(algo1_metrics).T
    df2 = pd.DataFrame(algo2_metrics).T

    # 综合评价
    score1 = topsis_evaluation(df1, weights).mean()
    score2 = topsis_evaluation(df2, weights).mean()

    subject_results[subject] = {
        '算法 1 得分': round(score1, 4),
        '算法 2 得分': round(score2, 4),
        '优势算法': '算法 1' if score1 > score2 else '算法 2',
        '样本量': len(data),
        '有效题型数': len(algo1_metrics)
    }

# 输出结果
print("\n=== 最终评价结果 ===")
result_df = pd.DataFrame(subject_results).T
print(result_df)

```

```

# 可视化展示
if not result_df.empty:
    plt.figure(figsize=(14, 6))

    # 学科对比图
    plt.subplot(121)
    result_df[['算法 1 得分', '算法 2 得分']].plot(kind='bar', ax=plt.gca())
    plt.title('各学科算法评分对比')
    plt.ylabel('综合得分')
    plt.xticks(rotation=45)
    plt.grid(axis='y', linestyle='--')

    # 样本量分布图
    plt.subplot(122)
    result_df['样本量'].plot(kind='pie', autopct='%1.1f%%', ax=plt.gca())
    plt.title('各学科样本量分布')
    plt.ylabel('')

    plt.tight_layout()
    plt.show()

if __name__ == "__main__":
    main()

```

8.4 问题四程序代码（基于误差阈值与分值构建 AI 算法动态分配与人工复核方案）

分配策略饼图

用途：反映各类题目采用不同算法的比例，便于直观展示资源分配。

```

import pandas as pd
import matplotlib.pyplot as plt

df = pd.read_excel('附件 3：科目题号、分值和误差阈值.xlsx')

# 假设分值 8 分以上且阈值 ≤ 1.5 分为 A，其余为 B
df['推荐算法'] = df.apply(lambda x: 'AI 算法 A+人工复核' if x['分值'] >= 8 and x['误差阈

```

```

值']<=1.5 else 'AI 算法 B+抽查', axis=1)
algo_counts = df['推荐算法'].value_counts()
plt.pie(algo_counts, labels=algo_counts.index, autopct='%1.1f%%')
plt.title('不同评分算法分配比例')
plt.show()

```

分值分布直方图

```

% 读取数据
data = readtable(filename, 'ReadVariableNames', false);
% 取分值列（第 2 列），并转换为数值型
score = data{:,2};
if iscell(score) || isstring(score)
    score = str2double(score);
end
% 去除可能的 NaN
score = score(~isnan(score));
% 绘制直方图
figure;
histogram(score, 'FaceColor', [0.2 0.6 1], 'EdgeColor', 'k');
xlabel('分值');
ylabel('题目数量');
title('分值分布直方图');
grid on;

```

AI 评分误差阈值分布箱型图

用途：分析 AI 评分误差范围，辅助判断阈值设置是否合理。

假设 df 中有一列为 AI 评分与人工评分误差

```

plt.boxplot(df['AI 评分误差'])
plt.title('AI 评分与人工评分误差分布')
plt.ylabel('误差分数')
plt.show()

```